

Framework for Multi-Cloud Using Integrity Verification Using CPDP Scheme

B. Shanmukhi*, D. Satyanarayana**

*[II-M. Tech] - CSE, Dr. KVSRRIT Kurnool.

**Asst. Prof in CSE Dept, Dr. KVSRRCEW Kurnool.

ABSTRACT

Provable data possession (PDP) is a technique for ensuring the integrity of data in storage outsourcing. In this paper, we address the construction of an efficient PDP scheme for distributed cloud storage to support the scalability of service and data migration, in which we consider the existence of multiple cloud service providers to cooperatively store and maintain the clients' data. We present a cooperative PDP (CPDP) scheme based on homomorphic verifiable response and hash index hierarchy. We prove the security of our scheme based on multi-prover zero-knowledge proof system, which can satisfy completeness, knowledge soundness, and zero-knowledge properties. In addition, we articulate performance optimization mechanisms for our scheme, and in particular present an efficient method for selecting optimal parameter values to minimize the computation costs of clients and storage service providers. Our experiments show that our solution introduces lower computation and communication overheads in comparison with non-cooperative approaches.

Index Terms—Storage Security, Provable Data Possession, Interactive Protocol, Zero-knowledge, Multiple Cloud, Cooperative

I. INTRODUCTION

In recent years, cloud storage service has become a faster profit growth point by providing a comparably low-cost, scalable, position-independent platform for clients' data. Since cloud computing environment is constructed based on open architectures and interfaces, it has the capability to incorporate multiple internal and/or external cloud services together to provide high interoperability. We call such a distributed cloud environment as a multi-Cloud (or hybrid cloud). Often, by using virtual infrastructure management (VIM), a multi-cloud allows clients to easily access his/her resources remotely through interfaces such as Web services provided by Amazon EC2. There exist various tools and technologies for multi-cloud, such as Platform VM Orchestrator, VMware vSphere, and Ovirt. These tools help cloud providers construct a distributed cloud storage platform (DCSP) for managing clients' data. However, if such an important platform is vulnerable to security attacks, it

would bring irretrievable losses to the clients. For example, the confidential data in an enterprise may be illegally accessed through a remote interface provided by a multi-cloud, or relevant data and archives may be lost or tampered with when they are stored into an uncertain storage pool outside the enterprise. Therefore, it is indispensable for cloud service providers (CSPs) to provide security techniques for managing their storage services.

Provable data possession (PDP) (or proofs of retrievability (POR)) is such a probabilistic proof technique for a storage provider to prove the integrity and ownership of clients' data without downloading data. The proof-checking without downloading makes it especially important for large-size files and folders (typically including many clients' files) to check whether these data have been tampered with or deleted without downloading the latest version of data. Thus, it is able to replace traditional hash and signature functions in storage outsourcing. Various PDP schemes have been recently proposed, such as Scalable PDP [4] and Dynamic PDP. However, these schemes mainly focus on PDP issues at untrusted servers in a single cloud storage provider and are not suitable for a multi-cloud environment (see the comparison of POR/PDP schemes in Table 1).

Motivation: To provide a low-cost, scalable, location independent platform for managing clients' data, current cloud storage systems adopt several new distributed file systems, for example, Apache Hadoop Distribution File System (HDFS), Google File System (GFS), Amazon S3 File System, CloudStore etc. These file systems share some similar features: a single metadata server provides centralized management by a global namespace; files are split into blocks or chunks and stored on block servers; and the systems are comprised of interconnected clusters of block servers. Those features enable cloud service providers to store and process large amounts of data. However, it is crucial to offer an efficient verification on the integrity and availability of stored data for detecting faults and automatic recovery. Moreover, this verification is necessary to provide reliability by automatically maintaining multiple copies of data and automatically redeploying processing logic in the event of failures.

Even though existing PDP schemes have addressed various security properties, such as public verifiability, dynamics, scalability, and privacy preservation, we still need a careful consideration of

some potential attacks, including two major categories: Data Leakage Attack by which an adversary can easily obtain the stored data through verification process after running or wiretapping sufficient verification communications (see Attacks 1 and 3 in Appendix A), and Tag Forgery Attack by which a dishonest CSP can deceive the clients (see Attacks 2 and 4 in Appendix A). These two attacks may cause potential risks for privacy leakage and ownership cheating. Also, these attacks can more easily compromise the security of a distributed cloud system than that of a single cloud system.

Although various security models have been proposed for existing PDP schemes, these models still cannot cover all security requirements, especially for provable secure privacy preservation and ownership authentication. To establish a highly effective security model, it is necessary to analyze the PDP scheme within the framework of zero-knowledge proof system (ZKPS) due to the reason that PDP system is essentially an interactive proof system (IPS), which has been well studied in the cryptography community. In summary, a verification scheme for data integrity in distributed storage environments should have the following features:

➤ Usability aspect: A client should utilize the integrity check in the way of collaboration services.

Table 1: Comparison of POR/PDP schemes for a file consisting of n blocks.

Scheme	CSP comp.	Client Comp.	Comm.	Frag.	Privacy	Dynamic Operations			Prob. of Detection
						modify	insert	delete	
PDP[5]	$O(t)$	$O(t)$	$O(1)$		✓				$1 - (1 - \rho)^t$
SPDP[6]	$O(t)$	$O(t)$	$O(t)$	✓	✓	✓ [#]		✓ [#]	$1 - (1 - \rho)^{t \cdot s}$
DPDP-I[7]	$O(t \log n)$	$O(t \log n)$	$O(t \log n)$		✓	✓	✓	✓	$1 - (1 - \rho)^t$
DPDP-II[7]	$O(t \log n)$	$O(t \log n)$	$O(t \log n)$			✓	✓	✓	$1 - (1 - \rho)^{\Omega(n)}$
CPOR-I[8]	$O(t)$	$O(t)$	$O(1)$						$1 - (1 - \rho)^t$
CPOR-II[8]	$O(t + s)$	$O(t + s)$	$O(s)$	✓					$1 - (1 - \rho)^{t \cdot s}$
Our Scheme	$O(t + s)$	$O(t + s)$	$O(s)$	✓	✓	✓	✓	✓	$1 - (1 - \rho)^{t \cdot s}$

s is the number of sectors in each block, c is the number of CSPs in a multi-cloud, t is the number of sampling blocks, ρ and ρk are the probability of block corruption in a cloud server and k -th cloud server in a multi-cloud $\mathcal{P} = \{Pk\}$, respectively, # denotes the verification process in a trivial approach, and $MHT, HomT, HomR$ denotes Merkle Hash tree, homomorphic tags, and homomorphic responses, respectively.

In order to support dynamic data operations, Ateniese et al. developed a dynamic PDP solution called Scalable PDP. They proposed a lightweight PDP scheme based on cryptographic hash function and symmetric key encryption, but the servers can deceive the owners by using previous metadata or responses due to the lack of randomness in the challenges. The numbers of updates and challenges are limited and fixed in advance and users cannot perform block insertions anywhere. Based on this work, Erway et al. introduced two Dynamic PDP schemes with a hash function tree to realize $(\log n)$

The scheme should conceal the details of the storage to reduce the burden on clients;

➤ Security aspect: The scheme should provide adequate security features to resist some existing attacks, such as data leakage attack and tag forgery attack;

➤ Performance aspect: The scheme should have the lower communication and computation overheads than non-cooperative solution.

➤ Related Works: To check the availability and integrity of outsourced data in cloud storages, researchers have proposed two basic approaches called Provable Data Possession (PDP) and Proofs of Retrievability (POR). Ateniese et al. first proposed the PDP model for ensuring possession of files on untrusted storages and provided an RSA-based scheme for a static case that achieves the (1) communication cost. They also proposed a publicly verifiable version, which allows anyone, not just the owner, to challenge the server for data possession. This property greatly extended application areas of PDP protocol due to the separation of data owners and the users. Moreover, they do not fit for multi-cloud storage due to the loss of homomorphism property in the verification process.

communication and computational costs for a n -block file. The basic scheme, called DPDP-I, retains the drawback of Scalable PDP, and in the ‘blockless’ scheme, called DPDP-II, the data blocks $\{m_{ij} \mid j \in [1, t]\}$ can be leaked by the response of a challenge, $M = \sum_{j=1}^t a_j m_{ij}$, where a_j is a random challenge value. Furthermore, these schemes are also not effective for a multi-cloud environment because the verification path of the challenge block cannot be stored completely in a cloud.

Juels and Kaliski [3] presented a POR scheme, which relies largely on preprocessing steps that the client conducts before sending a file to a CSP. Unfortunately, these operations prevent any efficient extension for updating data. Shacham and Waters proposed an improved version of this protocol called Compact POR, which uses homomorphic property to aggregate a proof into (1) authenticator value and $O(t)$ computation cost for t challenge blocks, but their solution is also static and could not prevent the leakage of data blocks in the verification process. Wang et al. presented a dynamic scheme

with $(\log n)$ cost by integrating the Compact POR scheme and Merkle Hash Tree (MHT) into the DPDP. Furthermore, several POR schemes and models have been recently proposed including [9], [10]. In [9] Bowers *et al.* introduced a distributed cryptographic system that allows a set of servers to solve the PDP problem.

This system is based on an integrity-protected error correcting code (IP-ECC), which improves the security and efficiency of existing tools, like POR. However, a file must be transformed into l distinct segments with the same length, which are distributed across l servers. Hence, this system is more suitable for RAID rather than cloud storage.

Our Contributions: In this paper, we address the problem of provable data possession in distributed cloud environments from the following aspects: *high security, transparent verification, and high performance*. To achieve these goals, we first propose a verification framework for multi-cloud storage along with two fundamental techniques: hash index hierarchy (HIH) and homomorphic verifiable response (HVR).

We then demonstrate that the possibility of constructing a cooperative PDP (CPDP) scheme without compromising data privacy based on modern cryptographic techniques, such as interactive proof system (IPS). We further introduce an effective construction of CPDP scheme using above-mentioned structure. Moreover, we give a security analysis of our CPDP scheme from the IPS model. We prove that this construction is a multi-prover zero-knowledge proof system (MP-ZKPS), which has completeness, knowledge soundness, and zero-knowledge properties. These properties ensure that CPDP scheme can implement the security against *data leakage attack* and *tag forgery attack*. To improve the system performance with respect to our scheme, we analyze the performance of probabilistic queries for detecting abnormal situations. This probabilistic method also has an inherent benefit in reducing computation and communication overheads. Then, we present an efficient method for the selection of optimal parameter values to minimize the computation overheads of CSPs and the clients' operations. In addition, we analyze that our scheme is suitable for existing distributed cloud storage systems. Finally, our experiments show that our solution introduces very limited computation and communication overheads.

Organization: The rest of this paper is organized as follows. In Section 2, we describe a formal definition of CPDP and the underlying techniques, which are utilized in the construction of our scheme. We introduce the details of cooperative PDP scheme for multicloud storage in Section 3. We describe the security and performance evaluation of our scheme in

Section 4 and 5, respectively. We discuss the related work in Section and Section 6 concludes this paper.

II. STRUCTURE AND TECHNIQUES

In this section, we present our verification framework for multi-cloud storage and a formal definition of CPDP. We introduce two fundamental techniques for constructing our CPDP scheme: hash index hierarchy (HIH) on which the responses of the clients' challenges computed from multiple CSPs can be combined into a single response as the final result; and homomorphic verifiable response (HVR) which supports distributed cloud storage in a multi-cloud storage and implements an efficient construction of collision-resistant hash function, which can be viewed as a random oracle model in the verification protocol.

2.1 Verification Framework for Multi-Cloud

Although existing PDP schemes offer a publicly accessible remote interface for checking and managing the tremendous amount of data, the majority of existing PDP schemes is incapable to satisfy the inherent requirements from multiple clouds in terms of communication and computation costs. To address this problem, we consider a multi-cloud storage service as illustrated in Figure 1. In this architecture, a data storage service involves three different entities: Clients who have a large amount of data to be stored in multiple clouds and have the permissions to access and manipulate stored data; Cloud Service Providers (CSPs) who work together to provide data storage services and have enough storages and computation resources; and Trusted Third Party (TTP) who is trusted to store verification parameters and offer public query services for these parameters.

In this architecture, we consider the existence of multiple CSPs to cooperatively store and maintain the clients' data. Moreover, a cooperative PDP is used to verify the integrity and availability of their stored data in all CSPs. The verification procedure is described as follows: Firstly, a client (data owner) uses the secret key to pre-process a file which consists of a collection of n blocks, generates a set of public verification information that is stored in TTP, transmits the file and some verification tags to CSPs, and may delete its local copy;

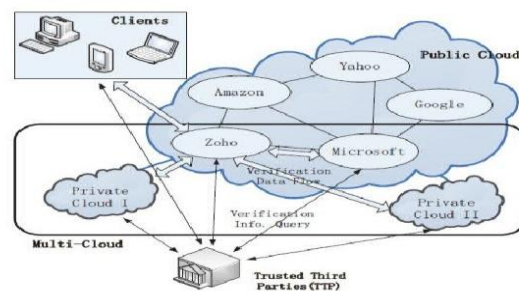


Fig 1: Verification architecture for data integrity

Then, by using a verification protocol, the clients can issue a challenge for one CSP to check the integrity and availability of outsourced data with respect to public information stored in TTP.

We neither assume that CSP is trust to guarantee the security of the stored data, nor assume that data owner has the ability to collect the evidence of the CSP's fault after errors have been found. To achieve this goal, a TTP server is constructed as a core trust base on the cloud for the sake of security. We assume the TTP is reliable and independent through the following functions: to setup and maintain the CPDP cryptosystem; to generate and store data owner's public key; and to store the public parameters used to execute the verification protocol in the CPDP scheme. Note that the TTP is not directly involved in the CPDP scheme in order to reduce the complexity of cryptosystem.

2.2 Definition of Cooperative PDP

In order to prove the integrity of data stored in a multi-cloud environment, we define a framework for CPDP based on interactive proof system (IPS) and multi-prover zero-knowledge proof system (MPZKPS), as follows:

Definition 1 (Cooperative-PDP): A cooperative provable data possession $\mathcal{S} = (KeyGen, TagGen, Proof)$ is a collection of two algorithms ($KeyGen, TagGen$) and an interactive proof system $Proof$, as follows: $K(1\kappa)$: takes a security parameter κ as input, and returns a secret key sk or a public-secret key pair (pk, sk) ;

$TagGen(sk, F, \mathcal{P})$: takes as inputs a secret key sk , a file F , and a set of cloud storage providers $\mathcal{P} = \{Pk\}$, and returns the triples (ζ, ψ, σ) , where ζ is the secret in tags, $\psi = (u, \mathcal{H})$ is a set of verification parameters u and an index hierarchy \mathcal{H} for F , $\sigma = \{\sigma(k)\} Pk \in \mathcal{P}$ denotes a set of all tags, $\sigma(k)$ is the tag of the fraction $F(k)$ of F in Pk ;

(\mathcal{P}, V) : is a protocol of proof of data possession between CSPs ($\mathcal{P} = \{Pk\}$) and a verifier (V), that is,

$$\left\langle \sum_{Pk \in \mathcal{P}} P_k(F^{(k)}, \sigma^{(k)}) \longleftrightarrow V \right\rangle (pk, \psi)$$

$$= \begin{cases} 1 & F = \{F^{(k)}\} \text{ is intact} \\ 0 & F = \{F^{(k)}\} \text{ is changed} \end{cases}$$

Where each Pk takes as input a file (k) and a set of tags (k) , and a public key pk and a set of public parameters ψ are the common input between P and V . At the end of the protocol run, V returns a bit $\{0|1\}$ denoting false and true. Where, $\sum_{Pk \in \mathcal{P}}$ denotes cooperative computing in $Pk \in \mathcal{P}$. A trivial way to realize the CPDP is to check the data stored in each cloud one by one, i.e.,

$$\bigwedge_{Pk \in \mathcal{P}} \langle P_k(F^{(k)}, \sigma^{(k)}) \longleftrightarrow V \rangle (pk, \psi),$$

Where \bigwedge denotes the logical AND operations among the boolean outputs of all protocols $\langle Pk, V \rangle$ for all $Pk \in \mathcal{P}$. However, it would cause significant communication and computation overheads for the verifier, as well as a loss of location-transparent. Such a primitive approach obviously diminishes the advantages of cloud storage: scaling arbitrarily up and down on demand. To solve this problem, we extend above definition by adding an organizer (O), which is one of CSPs that directly contacts with the verifier, as follows:

$$\left\langle \sum_{Pk \in \mathcal{P}} P_k(F^{(k)}, \sigma^{(k)}) \longleftrightarrow O \longleftrightarrow V \right\rangle (pk, \psi),$$

Where the action of organizer is to initiate and organize the verification process. This definition is consistent with aforementioned architecture, e.g., a client (or an authorized application) is considered as V , the CSPs are as $\mathcal{P} = \{Pi\} \in [1, c]$, and the Zoho cloud is as the organizer in Figure 1. Often, the organizer is an independent server or a certain CSP in \mathcal{P} . The advantage of this new multi-prover proof system is that it does not make any difference for the clients between multi-prover verification process and single-prover verification process in the way of collaboration. Also, this kind of transparent verification is able to conceal the details of data storage to reduce the burden on clients. For the sake of clarity, we list some used signals in Table 2.

Table 2: The signal and its explanation

Sig.	Repression
n	the number of blocks in a file;
s	the number of sectors in each block;
t	the number of index coefficient pairs in a query;
c	the number of clouds to store a file;
F	the file with $n \times s$ sectors, i.e., $F = \{mi, j\} i \in [1, n] j \in [1, s]$;
σ	the set of tags, i.e., $\sigma = \{\sigma i\} i \in [1, n]$;
Q	the set of index-coefficient pairs, i.e., $Q = \{(i, vi)\}$;
θ	The response for the challenge Q .

2.3 Hash Index Hierarchy for CPDP

To support distributed cloud storage, we illustrate a representative architecture used in our cooperative PDP scheme as shown in Figure 2. Our architecture has a hierarchy structure which resembles a natural representation of file storage. This hierarchical structure \mathcal{H} consists of three layers to represent relationships among all blocks for stored resources. They are described as follows:

- 1) **Express Layer:** offers an abstract representation of the stored resources;
- 2) **Service Layer:** offers and manages cloud storage services; and

3) **Storage Layer:** realizes data storage on many physical devices.

We make use of this simple hierarchy to organize data blocks from multiple CSP services into a large size file by shading their differences among these cloud storage systems. For example, in Figure 2 the resources in Express Layer are split and stored into three CSP's that are indicated by different colors, in Service Layer. In turn, each CSP fragments and stores the assigned data into the storage servers in Storage Layer. We also make use of colors to distinguish different CSPs. Moreover, we follow the logical order of the data blocks to organize the Storage Layer. This architecture also provides special functions for data storage and management, e.g., there may exist overlaps among data blocks (as shown in dashed boxes) and discontinuous blocks but these functions may increase the complexity of storage management.

In storage layer, we define a common fragment structure that provides probabilistic verification of data integrity for outsourced storage. The fragment structure is a data structure that maintains a set of block-tag pairs, allowing searches, checks and updates in (1) time. An instance of this structure is shown in storage layer of Figure 2: an outsourced file F is split into n blocks $\{m1,2, \dots, mn\}$, and each block mi is split into s sectors $\{mi,1,mi,2, \dots, mi,s\}$. The fragment structure consists of n block-tag pair $(mi, \sigma i)$, where σi is a signature tag of block mi generated by a set of secrets $\tau = (\tau1, \tau2, \dots, \tau s)$.

In order to check the data integrity, the fragment structure implements probabilistic verification as follows:

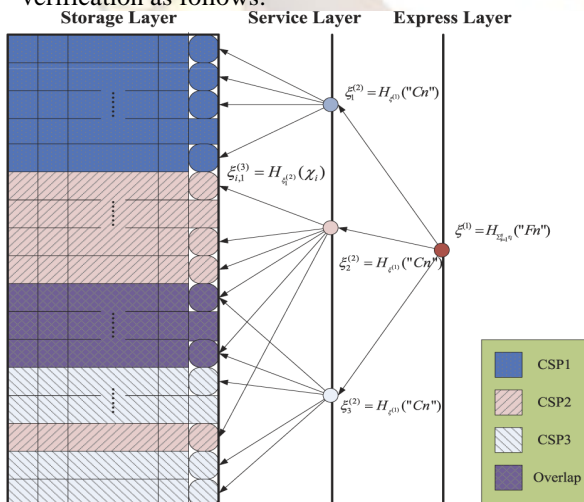


Fig 2. Index-hash hierarchy of CPDP model.

given a random chosen challenge (or query) $Q = \{(\cdot, vi)\}i \in I$, where I is a subset of the block indices and vi is a random coefficient. There exists an efficient algorithm to produce a constant-size response $(\mu1,$

$\mu2, \dots, \mu s, \sigma')$, where μi comes from all $\{mk,i, vk\}k \in I$ and σ' is from all $\{\sigma k, vk\}k \in I$.

Given a collision-resistant hash function (\cdot) , we make use of this architecture to construct a Hash Index Hierarchy \mathcal{H} (viewed as a random oracle), which is used to replace the common hash function in prior PDP schemes, as follows:

1) **Express layer:** given s random $\{\tau i\} i=1$ and the file name $F n$, sets $(1) = H \sum s i=1 (F n)$ and makes it public for verification but makes $\{\tau i\} i=1$ secret;

2) **Service layer:** given the (1) and the cloud name $C k$, sets $(2) k = (1) (C k)$;

3) **Storage layer:** given the (2), a block number i , and its index record $\chi i = "Bi||Vi||Ri"$, sets $(3) i, = (2)k (\chi i)$, where Bi is the sequence number of a block, Vi is the updated version number, and Ri is a random integer to avoid collision. As a virtualization approach, we introduce a simple index-hash table $\chi = \{\chi i\}$ to record the changes of file blocks as well as to generate the hash value of each block in the verification process. The structure of χ is similar to the structure of file block allocation table in file systems. The index-hash table consists of serial number, block number, version number, random integer, and so on. Different from the common index table, we assure that all records in our index table differ from one another to prevent forgery of data blocks and tags. By using this structure, especially the index records $\{\chi i\}$, our CPDP scheme can also support dynamic data operations. The proposed structure can be readily incorporated into MAC-based, ECC or RSA schemes. These schemes, built from collision-resistance signatures (see Section 3.1) and the random oracle model, have the shortest query and response with public verifiability. They share several common characters for the implementation of the CPDP framework in the multiple clouds:

1) a file is split into $n \times s$ sectors and each block (s sectors) corresponds to a tag, so that the storage of signature tags can be reduced by the increase of s ;

2) a verifier can verify the integrity of file in random sampling approach, which is of utmost importance for large files;

3) these schemes rely on homomorphic properties to aggregate data and tags into a constant size response, which minimizes the overhead of network communication; and

4) the hierarchy structure provides a virtualization approach to conceal the storage details of multiple CSPs.

1) a file is split into $n \times s$ sectors and each block (s sectors) corresponds to a tag, so that the storage of signature tags can be reduced by the increase of s ;

2) a verifier can verify the integrity of file in random sampling approach, which is of utmost importance for large files;

3) these schemes rely on homomorphic properties to aggregate data and tags into a constant size response, which minimizes the overhead of network communication; and

4) the hierarchy structure provides a virtualization approach to conceal the storage details of multiple CSPs.

2.4 Homomorphic Verifiable Response for CPDP

A homomorphism is a map: $\mathbb{P} \rightarrow \mathbb{Q}$ between two groups such that $(g1 \oplus g2) = (g1) \otimes (g2)$ for all $g1, g2 \in \mathbb{P}$, where \oplus denotes the operation in \mathbb{P} and \otimes denotes the operation in \mathbb{Q} . This notation has been used to define Homomorphic Verifiable Tags (HVTs) in [2]: Given two values σi and σj for two messages mi and anyone can combine them into a value σ' corresponding to the sum of the messages mi

+ m_j . When provable data possession is considered as a challenge-response protocol, we extend this notation to the concept of Homomorphic Verifiable Responses (HVR), which is used to integrate multiple responses from the different CSPs in CPDP scheme as follows:

Definition 2 (Homomorphic Verifiable Response): A response is called homomorphic verifiable response in a PDP protocol, if given two responses θ_i and θ_j for two challenges Q_i and Q_j from two CSPs, there exists an efficient algorithm to combine them into a response θ corresponding to the sum of the challenges $Q_i \cup Q_j$. Homomorphic verifiable response is the key technique of CPDP because it not only reduces the communication bandwidth, but also conceals the location of outsourced data in the distributed cloud storage environment.

III. COOPERATIVE PDP SCHEME

In this section, we propose a CPDP scheme for multicloud system based on the above-mentioned structure and techniques. This scheme is constructed on collision-resistant hash, bilinear map group, aggregation algorithm, and homomorphic responses.

3.1 Notations and Preliminaries

Let $\mathbb{H} = \{Hk\}$ be a family of hash functions: $\{0, 1\}^n \rightarrow \{0, 1\}^*$ index by $k \in \mathcal{K}$. We say that algorithm \mathcal{A} has advantage ϵ in breaking collision resistance of \mathbb{H} if $\Pr[\mathcal{A}(k) = (m_0, m_1) : m_0 \neq m_1, Hk(m_0) = Hk(m_1)] \geq \epsilon$, where the probability is over the random choices of $k \in \mathcal{K}$ and the random bits of \mathcal{A} . So that, we have the following definition.

Definition 3 (Collision-Resistant Hash): A hash family \mathbb{H} is (t, ϵ) -collision-resistant if no t -time adversary has advantage at least ϵ in breaking collision resistance of \mathbb{H} . We set up our system using bilinear pairings proposed by Boneh and Franklin [14]. Let \mathbb{G} and $\mathbb{G}T$ be two multiplicative groups using elliptic curve conventions with a large prime order p . The function e is a computable bilinear map $e : \mathbb{G} \times \mathbb{G} \rightarrow \mathbb{G}T$ with the following properties: for any $G, H \in \mathbb{G}$ and all $a, b \in \mathbb{Z}p$, we have 1) Bilinearity: $e([a]G, [b]H) = e(G, H)ab$; 2) Non-degeneracy: $e(G, H) \neq 1$ unless G or $H = 1$; and 3) Computability: (G, \cdot) is efficiently computable.

Definition 4 (Bilinear Map Group System): A bilinear map group system is a tuple $\mathbb{S} = \langle p, \cdot, e \rangle$ composed of the objects as described above.

3.2 Our CPDP Scheme

In our scheme (see Fig 3), the manager first runs algorithm *KeyGen* to obtain the public/private key pairs for CSPs and users. Then, the clients generate the tags of outsourced data by using *TagGen*. Anytime, the protocol *Proof* is performed by a 5-move interactive proof protocol between a verifier and more than one CSP, in which CSPs need

not to interact with each other during the verification process, but an organizer is used to organize and manage all CSPs.

This protocol can be described as follows: 1) the organizer initiates the protocol and sends a commitment to the verifier; 2) the verifier returns a challenge set of random index-coefficient pairs Q to the organizer; 3) the organizer relays them into each P_i in \mathcal{P} according to the exact position of each data block; 4) each P_i returns its response of challenge to the organizer; and 5) the organizer synthesizes a final response from received responses and sends it to the verifier.

The above process would guarantee that the verifier accesses files without knowing on which CSPs or in what geographical locations their files reside. In contrast to a single CSP environment, our scheme differs from the common PDP scheme in two aspects:

1) Tag aggregation algorithm: In stage of commitment, the organizer generates a random $\gamma \in \mathbb{R} \mathbb{Z}p$ and returns its commitment H_γ to the verifier. This assures that the verifier and CSPs do not obtain the value of γ . Therefore, our approach guarantees only the organizer can compute the final σ' by using γ and σ'_k received from CSPs.

After σ' is computed, we need to transfer it to the organizer in stage of "Response1". In order to ensure the security of transmission of data tags, our scheme employs a new method, similar to the E_1 Gamal encryption, to encrypt the combination of tags $\pi_{(i,v_i) \in Q} k \sigma^{v_i}$, that is, for $sk = s \in \mathbb{Z}p$ and $pk = (g, S = g^s) \in \mathbb{G}_2$, the cipher of message m is $C = (C1 = gr, C2 = m \cdot S^r)$ and its decryption is performed by $m = C2 \cdot C^{-s_1}$. Thus, we hold the equation

$$\begin{aligned} \sigma' &= \left(\prod_{P_k \in \mathcal{P}} \frac{\sigma'_k}{\eta_k^s} \right)^\gamma = \left(\prod_{P_k \in \mathcal{P}} \frac{S^{r_k} \cdot \prod_{(i,v_i) \in Q_k} \sigma_i^{v_i}}{\eta_k^s} \right)^\gamma \\ &= \left(\prod_{P_k \in \mathcal{P}} \cdot \prod_{(i,v_i) \in Q_k} \sigma_i^{v_i} \right)^\gamma = \prod_{(i,v_i) \in Q} \sigma_i^{v_i \cdot \gamma}. \end{aligned}$$

2) Homomorphic responses: Because of the homomorphic property, the responses computed from CSPs in a multi-cloud can be combined into a single final response as follows: given a set of $\theta_k = (\pi_k, \sigma'_k, \mu_k, \eta_k)$ received from P_k , let $\lambda_j = \sum_{P_k \in \mathcal{P}} \lambda_{j,k}$, the organizer can compute

$$\begin{aligned} \mu'_j &= \sum_{P_k \in \mathcal{P}} \gamma \cdot \mu_{j,k} = \sum_{P_k \in \mathcal{P}} \gamma \cdot \left(\lambda_{j,k} + \sum_{(i,v_i) \in Q_k} v_i \cdot m_{i,j} \right) \\ &= \sum_{P_k \in \mathcal{P}} \gamma \cdot \lambda_{j,k} + \gamma \cdot \sum_{P_k \in \mathcal{P}} \sum_{(i,v_i) \in Q_k} v_i \cdot m_{i,j} \\ &= \gamma \cdot \sum_{P_k \in \mathcal{P}} \lambda_{j,k} + \gamma \cdot \sum_{(i,v_i) \in Q} v_i \cdot m_{i,j} \\ &= \gamma \cdot \lambda_j + \gamma \cdot \sum_{(i,v_i) \in Q} v_i \cdot m_{i,j}. \end{aligned}$$

The commitment of λ_j is also computed by

$$\begin{aligned} \pi^j &= \left(\prod_{P_k \in \mathcal{P}} \pi_k \right)^\gamma = \left(\prod_{P_k \in \mathcal{P}} \prod_{j=1}^s \pi_{j,k} \right)^\gamma \\ &= \prod_{j=1}^s \prod_{P_k \in \mathcal{P}} e(u_j^{\lambda_{j,k}}, H_2)^\gamma \\ &= \prod_{j=1}^s e(u_j^{\sum_{P_k \in \mathcal{P}} \lambda_{j,k}}, H_2)^\gamma = \prod_{j=1}^s e(u_j^{\lambda_j}, H_2^j). \end{aligned}$$

It is obvious that the final response θ received by the verifier's from multiple CSPs is same as that in one simple CSP. This means that our CPDP scheme is able to provide a transparent verification for the verifiers. Two response algorithms, Response1 and Response2, comprise an HVR: Given two responses θ_i and θ_j for two challenges Q_i and Q_j from two CSPs, i.e., $\theta_i = \text{Response1}(Q_i, \{m_k\}_{k \in I_i}, \{\sigma_k\}_{k \in I_i})$, there exists an efficient algorithm to combine them into a final response θ corresponding to the sum of the challenges

$$\begin{aligned} Q_i \cup Q_j, \text{ that is,} \\ \theta &= \text{Response1}(Q_i \cup Q_j, \{m_k\}_{k \in I_i \cup I_j}, \{\sigma_k\}_{k \in I_i \cup I_j}) \\ &= \text{Response2}(\theta_i, \theta_j). \end{aligned}$$

For multiple CSPs, the above equation can be extended to $\theta = \text{Response2}(\{\theta_k\}_{k \in \mathcal{P}})$. More importantly, the HVR is a pair of values $\theta = (\pi, \sigma, \mu)$, which has a constant-size even for different challenges.

IV. SECURITY ANALYSIS

We give a brief security analysis of our CPDP construction. This construction is directly derived from multi-prover zero-knowledge proof system (MPZKPS), which satisfies following properties for a given assertion, L :

- 1) **Completeness:** whenever $x \in L$, there exists a strategy for the provers that convinces the verifier that this is the case;
- 2) **Soundness:** whenever $x \notin L$, whatever strategy the provers employ, they will not convince the verifier that $x \in L$;
- 3) **Zero-knowledge:** no cheating verifier can learn anything other than the veracity of the statement.

According to existing IPS research, these properties can protect our construction from various attacks, such as data leakage attack (privacy leakage), tag forgery attack (ownership cheating), etc. In details, the security of our scheme can be analyzed as follows:

4.1 Collision resistant for index-hash hierarchy

In our CPDP scheme, the collision resistant of index hash hierarchy is the basis and prerequisite for the security of whole scheme, which is described as being secure in the *random oracle model*. Although the hash function is collision resistant, a successful hash collision can still be used to produce a forged tag when the same hash value is reused multiple times, e.g., a legitimate client modifies the data or repeats to insert and delete data blocks of outsourced data. To avoid the hash collision, the hash value $\xi(3)_{i,k}$, which is used to generate the tag σ_i in

CPDP scheme, is computed from the set of values $\{\tau_i\}, Fn, Ck, \{\chi_i\}$. As long as there exists one bit difference in these data, we can avoid the hash collision. As a consequence, we have the following theorem (see Appendix B):

Theorem 1 (Collision Resistant): The index-hash hierarchy in CPDP scheme is collision resistant, even if the client generates $\sqrt{2p} \cdot \ln(1/(1-\epsilon))$ files with the same file name and cloud name, and the client repeats $\sqrt{2L+1} \cdot \ln(1/(1-\epsilon))$ times to modify, insert and delete data blocks, where the collision probability is at least ϵ , $\tau_i \in \mathbb{Z}_p$, and $|Ri| = L$ for $Ri \in \chi_i$.

4.2 Completeness property of verification

In our scheme, the completeness property implies public verifiability property, which allows anyone, not just the client (data owner), to challenge the cloud server for *data integrity* and *data ownership* without the need for any secret information. First, for every available data-tag pair $(F, \sigma) \in Ta(sk, F)$ and a random challenge $Q = (i, v_i)_{i \in I}$, the verification protocol should be completed with success probability according to the Equation (3), that is,

$$\Pr \left[\left\langle \sum_{P_k \in \mathcal{P}} F_k(F^{(k)}, \sigma^{(k)}) \leftrightarrow O \leftrightarrow V \right\rangle (pk, \psi) = 1 \right] = 1.$$

In this process, anyone can obtain the owner's public key $pk = (g, h, 1 = h\alpha, H_2 = h\beta)$ and the corresponding file parameter $\psi = (u, \xi(1), \chi)$ from TTP to execute the verification protocol, hence this is a public verifiable protocol. Moreover, for different owners, the secrets α and β hidden in their public key pk are also different, determining that a success verification can only be implemented by the real owner's public key. In addition, the parameter ψ is used to store the file-related information, so an owner can employ a unique public key to deal with a large number of outsourced files.

4.3 Zero-knowledge property of verification

The CPDP construction is in essence a Multi-Prover Zero-knowledge Proof (MP-ZKP) system, which can be considered as an extension of the notion of

$$\begin{aligned} \pi^j \cdot e(\sigma^j, h) &= \prod_{j=1}^s e(u_j^{\lambda_j}, H_2) \cdot e\left(\prod_{(i,v_i) \in Q} \sigma_i^{v_i}, h\right) \\ &= \prod_{j=1}^s e(u_j^{\lambda_j}, H_2) \cdot e\left(\prod_{(i,k) \in Q} (\xi_{i,k}^{(3)})^\alpha \cdot \left(\prod_{j=1}^s u_j^{m_{i,j}}\right)^\beta, h\right) \\ &= \prod_{j=1}^s e(u_j^{\lambda_j}, H_2) \cdot e\left(\prod_{(i,v_i) \in Q} (\xi_{i,k}^{(3)})^{v_i} \cdot h\right)^{\alpha\gamma} \cdot e\left(\prod_{j=1}^s u_j^{\sum_{(i,v_i) \in Q} v_i m_{i,j}}, h^\beta\right) \\ &= e\left(\prod_{(i,v_i) \in Q} (\xi_{i,k}^{(3)})^{v_i}, H_1\right) \cdot \prod_{j=1}^s e(u_j^{\lambda_j}, H_2). \end{aligned}$$

an interactive proof system (IPS). Roughly speaking, in the scenario of MP-ZKP, a polynomial-time bounded verifier interacts with several provers whose computational powers are unlimited. According to a *Simulator* model, in which every cheating verifier has a simulator that can produce a transcript that "looks like" an interaction between a honest prover and a

cheating verifier, we can prove our CPDP construction has Zero-knowledge property (see Appendix C):

Theorem 2 (Zero-Knowledge Property): The verification protocol $Proof(\mathcal{P}, V)$ in CPDP scheme is a computational zero-knowledge system under a simulator model, that is, for every probabilistic polynomial-time interactive machine V_+ , there exists a probabilistic polynomial-time algorithm S^* such that the ensembles $V_{i ew}(\langle \sum_{Pk \in \mathcal{P}} Pk(F(k), \sigma(k)) \leftrightarrow O \leftrightarrow V^* \rangle (pk, \psi))$ and $S^*(pk, \psi)$ are computationally indistinguishable. Zero-knowledge is a property that achieves the CSPs' robustness against attempts to gain knowledge by interacting with them. For our construction, we make use of the zero-knowledge property to preserve the privacy of data blocks and signature tags. Firstly, randomness is adopted into the CSPs' responses in order to resist the *data leakage attacks* (see Attacks 1 and 3 in Appendix A). That is, the random integer $\lambda_{j,k}$ is introduced into the response $\mu_{j,k}$, i.e., $\mu_{j,\square} = \square_{\square,\square} + \sum(\square, \square) \in \square \square \square \square \cdot \square \square, \square$. This means that the cheating verifier cannot obtain $\square \square, \square$ from $\square \square, \square$ because he does not know the random integer $\square \square, \square$. At the same time, a random integer \square is also introduced to randomize the verification tag \square , i.e., $\square' \leftarrow (\prod \square \square \in \square \square' \square \cdot \square - \square \square) \square$. Thus, the tag \square cannot reveal to the cheating verifier in terms of randomness.

4.4 Knowledge soundness of verification

For every data-tag pairs $(\square^*, \square^*) \notin \square \square \square \square \square \square (\square \square, \square)$, in order to prove nonexistence of fraudulent \square^* and \square^* , we require that the scheme satisfies the knowledge soundness property, that is,

$$\Pr \left[\left\langle \sum_{Pk \in \mathcal{P}^*} Pk(F^{(k)*}, \sigma^{(k)*}) \leftrightarrow O^* \leftrightarrow V \right\rangle (pk, \psi) = 1 \right] \leq \epsilon,$$

Where ϵ is a negligible error. We prove that our scheme has the knowledge soundness property by using reduction to absurdity 1: we make use of \square^* to construct a knowledge extractor \mathcal{M} which gets the common input $(\square \square, \square)$ and rewindable black box accesses to the prover \square^* , and then attempts to break the computational Diffie-Hellman (CDH) problem in \square : given $\square, \square^1 = \square \square, \square^2 = \square \square \in \square \square$, output $\square^{\square \square} \in \square$. But it is unacceptable because the CDH problem is widely regarded as an unsolved problem in polynomial-time. Thus, the opposite direction of the theorem also follows.

Theorem 3 (Knowledge Soundness Property): Our scheme has (\square, \square') knowledge soundness in random oracle and rewind able knowledge extractor model assuming the (\square, \square) -computational Diffie-Hellman (CDH) assumption holds in the group \square for $\square' \geq \square$. Essentially, the soundness means that it is infeasible to fool the verifier to accept false statements. Often,

the soundness can also be regarded as a stricter notion of unforgeability for file tags to avoid cheating the ownership. This means that the CSPs, even if collusion is attempted, cannot be tampered with the data or forge the data tags if the soundness property holds. Thus, the Theorem 3 denotes that the CPDP scheme can resist the *tag forgery attacks* to avoid cheating the CSPs' ownership.

V. PERFORMANCE EVALUATIONS

In this section, to detect abnormality in a low overhead and timely manner, we analyze and optimize the performance of CPDP scheme based on the above scheme from two aspects: evaluation of probabilistic queries and optimization of length of blocks. To validate the effects of scheme, we introduce a prototype of CPDP-based audit system and present the experimental results.

5.1 Performance Analysis for CPDP Scheme

We present the computation cost of our CPDP scheme in Table 3. We use $[\square]$ to denote the computation cost of an exponent operation in \square , namely, $\square \square$, where \square is a positive integer in $\mathbb{Z} \square$ and $\square \in \square$ or $\square \square$. We neglect the computation cost of algebraic operations and simple modular arithmetic operations because they run fast enough [16]. The most complex operation is the computation of a bilinear map (\cdot, \cdot) between two elliptic points (denoted as $[\square]$).

Table 3: Comparison of computation overheads between our CPDP scheme and non-cooperative (trivial) scheme.

	CPDP Scheme	Trivial Scheme
KeyGen	3 $[\square]$	2 $[E]$
TagGen	(2 $\square + s$) $[\square]$	(2 $\square + \square$) $[\square]$
Proof(\square)	$\square[\square] + (\square + \square \square + 1)[\square]$	$\square[\square] + (\square + \square \square - \square)[\square]$
Proof(V)	3 $[\square] + (\square + \square)[\square]$	3 $\square[\square] + (\square + \square \square)[\square]$

Then, we analyze the storage and communication costs of our scheme. We define the bilinear pairing takes the form: $\square(\square \square \square) \times \square(\square \square \square) \rightarrow \square^* \square \square \square$ (The definition given here is from [17], [18]), where \square is a prime, \square is a positive integer, and \square is the embedding degree (or security multiplier). In this case, we utilize an asymmetric pairing: $\square^1 \times \square^2 \rightarrow \square \square$ to replace the symmetric pairing in the original schemes. In Table 3, it is easy to find that client's computation overheads are entirely irrelevant for the number of CSPs. Further, our scheme has better performance compared with non-cooperative approach due to the total of computation overheads decrease $3(\square - 1)$ times bilinear map operations, where \square is the number of clouds in a multicloud. The reason is that, before the responses are sent to the verifier from \square clouds, the organizer has aggregate these responses into a

response by using aggregation algorithm, so the verifier only need to verify this response once to obtain the final result.

Table 4: Comparison of communication overheads between our CPDP and non-cooperative (trivial) scheme.

	CPDP Scheme	Trivial Scheme
Commitment	κ_2	cl_2
Challenge 1	$2tl_0$	$2tl_0$
Challenge 2	$2tl_0/c$	
Response1	$sl_0 + 2l_I + l_T$	$(sl_0 + l_I + l_T)c$
Response2	$sl_0 + l_I + l_T$	

Without loss of generality, let the security parameter κ be 80 bits, we need the elliptic curve domain parameters over \mathbb{F}_q with $|\kappa| = 160$ bits and $\kappa = 1$ in our experiments. This means that the length of integer is $\kappa_0 = 2\kappa$ in \mathbb{Z} . Similarly, we have $\kappa_1 = 4\kappa$ in \mathbb{F}_1 , $\kappa_2 = 24\kappa$ in \mathbb{F}_2 , and $\kappa_3 = 24\kappa$ in \mathbb{F}_3 for the embedding degree $\kappa = 6$. The storage and communication cost of our scheme is shown in Table 4. The storage overhead of a file with $\kappa(\kappa) = 1\kappa$ -bytes is $\kappa(\kappa) = \kappa \cdot \kappa \cdot \kappa_0 + \kappa \cdot \kappa_1 = 1.04\kappa$ -bytes for $\kappa = 103$ and $\kappa = 50$. The storage overhead of its index table κ is $\kappa \cdot \kappa_0 = 20\kappa$ -bytes. We define the overhead rate as $\kappa = \kappa(\kappa) / \kappa(\kappa) - 1 = \kappa_1 / \kappa_0$ and it should therefore be kept as low as possible in order to minimize the storage in cloud storage providers. It is obvious that a higher κ means much lower storage. Furthermore, in the verification protocol, the communication overhead of challenge is $2\kappa \cdot \kappa_0 = 40 \cdot \kappa$ -Bytes in terms of the number of challenged blocks κ , but its response (response1 or response2) has a constant-size communication overhead $\kappa \cdot \kappa_0 + \kappa_1 + \kappa_3 \approx 1.3\kappa$ -bytes for different file sizes. Also, it implies that client's communication overheads are of a fixed size, which is entirely irrelevant for the number of CSPs.

5.2 Parameter Optimization

In the fragment structure, the number of sectors per block κ is an important parameter to affect the performance of storage services and audit services. Hence, we propose an optimization algorithm for the value of s in this section. Our results show that the optimal value can not only minimize the computation and communication overheads, but also reduce the size of extra storage, which is required to store the verification tags in CSPs. Assume κ denotes the probability of sector corruption. In the fragment structure, the choosing of κ is extremely important for improving the performance of the CPDP scheme. Given the detection probability κ and the probability of sector corruption ρ for multiple clouds $\kappa = \{\kappa_1, \kappa_2, \kappa_3\}$, the optimal value of κ can be computed by $\min_{s \in \mathbb{N}} \{ \log(1-\kappa) / (\sum_{\rho_k \in \mathcal{P}} \rho_k \cdot \log(1-\rho_k)) \cdot \kappa / \kappa + \kappa \cdot \kappa + \kappa \}$, where $\kappa \cdot \kappa + \kappa \cdot \kappa + \kappa$ denotes the computational cost of verification protocol in PDP scheme, $\kappa, \kappa, \kappa \in \mathbb{R}$, and κ is a constant. This

conclusion can be obtained from following process: Let $\kappa = \kappa \cdot \kappa = \kappa(\kappa) / l_0$. According to above-mentioned results, the sampling probability holds $\kappa \geq (\log(1-\kappa)) / (\kappa \cdot \sum_{\rho_k \in \mathcal{P}} \rho_k \cdot \log(1-\rho_k)) = (\log(1-\kappa)) / \kappa \cdot \sum_{\rho_k \in \mathcal{P}} \rho_k \cdot \log(1-\rho_k)$. In order to minimize the computational cost, we have

$$\begin{aligned} & \min_{s \in \mathbb{N}} \{ a \cdot t + b \cdot s + c \} \\ & = \min_{s \in \mathbb{N}} \{ a \cdot n \cdot w + b \cdot s + c \} \\ & \geq \min_{s \in \mathbb{N}} \left\{ \frac{\log(1-P)}{\sum_{\rho_k \in \mathcal{P}} \rho_k \cdot \log(1-\rho_k)} \frac{a}{s} + b \cdot s + c \right\}. \end{aligned}$$

where κ denotes the proportion of data blocks in the κ -th CSP, ρ_k denotes the probability of file corruption in the κ -th CSP. Since κ / κ is a monotone decreasing function and $\kappa \cdot \kappa$ is a monotone increasing function for $\kappa > 0$, there exists an optimal value of $\kappa \in \mathbb{N}$ in the above equation. The optimal value of κ is unrelated to a certain file from this conclusion if the probability ρ is a constant value. For instance, we assume a multi-cloud storage involves three CSPs $\kappa = \{\kappa_1, \kappa_2, \kappa_3\}$ and the probability of sector corruption is a constant value $\{\rho_1, \rho_2, \rho_3\} = \{0.01, 0.02, 0.001\}$. We set the detection probability κ with the range from 0.8 to 1, e.g., $\{0.8, 0.85, 0.9, 0.95, 0.99, 0.999\}$. For a file, the proportion of data blocks is 50%, 30%, and 20% in three CSPs, respectively, that is, $\kappa_1 = 0.5, \kappa_2 = 0.3, \text{ and } \kappa_3 = 0.2$. In terms of Table 3, the computational cost of CSPs can be simplified to $\kappa + 3\kappa + 9$. Then, we can observe the computational cost under different κ and κ in Figure 4. When κ is less than the optimal value, the computational cost decreases evidently with the increase of κ , and then it raises when κ is more than the optimal value.

More accurately, we show the influence of parameters, $\kappa \cdot \kappa, \kappa,$ and $\kappa,$ under different detection probabilities in Table 6. It is easy to see that computational cost rises with the increase of κ . Moreover, we can make sure the sampling number of challenge with following conclusion:

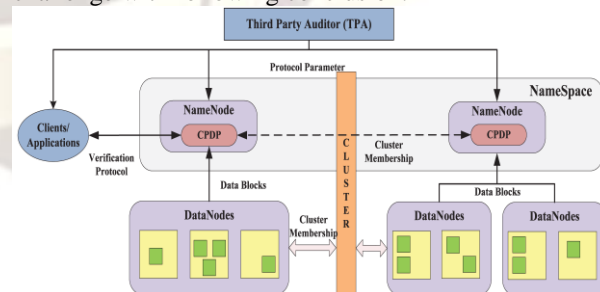


Fig 3: Applying CPDP scheme in Hadoop distributed file system (HDFS)

Given the detection probability κ , the probability of sector corruption ρ , and the number of sectors in each block κ , the sampling number of verification protocol are a constant $\kappa = \kappa \cdot \kappa$.

$\frac{1}{\alpha} \geq (\log(1-\alpha)) / \alpha \cdot \sum_{i=1}^n \alpha_i \cdot \log(1-\alpha_i)$ for different files.

Table 6: The influence of parameters under different detection probabilities α ($\alpha = \{\alpha_1, \alpha_2, \alpha_3\} = \{0.01, 0.02, 0.001\}$, $\{\alpha_1, \alpha_2, \alpha_3\} = \{0.5, 0.3, 0.2\}$).

P	0.8	0.85	0.9	0.95	0.99	0.999
S	142.	168.0	204.0	265.4	408.0	612.0
z	6	9	2	3	4	6
w						
s	7	8	10	11	13	16
w	20	21	20	29	31	38

Finally, we observe the change of α under different α and β . The experimental results are shown in Table 5. It is obvious that the optimal value of α rises with increase of β and with the decrease of β . We choose the optimal value of α on the basis of practical settings and system requisition. For NTFS format, we suggest that the value of α is 200 and the size of block is 4KBytes, which is the same as the default size of cluster when the file size is less than 16TB in NTFS. In this case, the value of α ensures that the extra storage doesn't exceed 1% in storage servers.

5.3 CPDP for Integrity Audit Services

Based on our CPDP scheme, we introduce audit system architecture for outsourced data in multiple clouds by replacing the TTP with a third party auditor (TPA) in Figure 1. In this architecture, this architecture can be constructed into a visualization infrastructure of cloud-based storage service. In Figure 5, we show an example of applying our CPDP scheme in Hadoop distributed file system (HDFS) 4, which a distributed, scalable, and portable file system. HDFS' architecture is composed of NameNode and DataNode, where NameNode maps a file name to a set of indexes of blocks and DataNode indeed stores data blocks. To support our CPDP scheme, the index-hash hierarchy and the metadata of NameNode should be integrated together to provide an enquiry service for the hash value $(h_i)_{i=1}^n$ or index-hash record $(i)_{i=1}^n$. Hence, it is easy to replace the checksum methods with the CPDP scheme for anomaly detection in current HDFS. To validate the effectiveness and efficiency of our proposed approach for audit services, we have implemented a prototype of an audit system. We simulated the audit service and the storage service by using two local IBM servers with two Intel Core 2 processors at 2.16 GHz and 500M RAM running Windows Server 2003. These servers were connected via 250 MB/sec of network bandwidth. Using GMP and PBC libraries, we have implemented a cryptographic library upon which our scheme can be constructed. This C library

contains approximately 5,200 lines of codes and has been tested on both Windows and Linux platforms. The elliptic curve utilized in the experiment is a MNT curve, with base field size of 160 bits and the embedding degree 6. The security level is chosen to be 80 bits, which means $|G| = 160$.

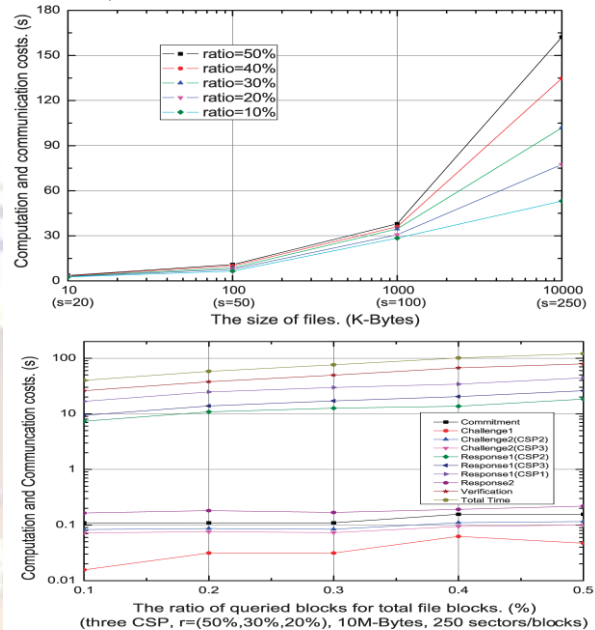


Fig 4: Experimental results under different file size, sampling ratio, and sector number

Firstly, we quantify the performance of our audit scheme under different parameters, such as file size β , sampling ratio α , sector number per block β , and so on. Our analysis shows that the value of α should grow with the increase of β in order to reduce computation and communication costs. Thus, our experiments were carried out as follows: the stored files were chosen from 10KB to 10MB; the sector numbers were changed from 20 to 250 in terms of file sizes; and the sampling ratios were changed from 10% to 50%. The experimental results are shown in the left side of Figure 6. These results dictate that the computation and communication costs (including I/O costs) grow with the increase of file size and sampling ratio. Next, we compare the performance of each activity in our verification protocol. We have shown the theoretical results in Table 4: the overheads of "commitment" and "challenge" resemble one another, and the overheads of "response" and "verification" resemble one another as well. To validate the theoretical results, we changed the sampling ratio α from 10% to 50% for a 10MB file and 250 sectors per block in a multi-cloud $\alpha = \{\alpha_1, \alpha_2, \alpha_3\}$, in which the proportions of data blocks are 50%, 30%, and 20% in three CSPs, respectively. In the right side of Figure 6, our experimental results show that the computation and communication costs of "commitment" and "challenge" are slightly changed along with the sampling ratio, but those for "response" and "verification" grows with the increase of the

sampling ratio. Here, “challenge” and “response” can be divided into two sub-processes: “challenge1” and “challenge2”, as well as “response1” and “response2”, respectively. Furthermore, the proportions of data blocks in each CSP have greater influence on the computation costs of “challenge” and “response” processes. In summary, our scheme has better performance than non-cooperative approach.

VI. CONCLUSIONS

In this paper, we presented the construction of an efficient PDP scheme for distributed cloud storage. Based on homomorphic verifiable response and hash index hierarchy, we have proposed a cooperative PDP scheme to support dynamic scalability on multiple storage servers. We also showed that our scheme provided all security properties required by zero knowledge interactive proof system, so that it can resist various attacks even if it is deployed as a public audit service in clouds. Furthermore, we optimized the probabilistic query and periodic verification to improve the audit performance. Our experiments clearly demonstrated that our approaches only introduce a small amount of computation and communication overheads. Therefore, our solution can be treated as a new candidate for data integrity verification in outsourcing data storage systems. Finally, it is still a challenging problem for the generation of tags with the length irrelevant to the size of data blocks. We would explore such an issue to provide the support of variable-length block verification.

REFERENCES

- [1]. B. Sotomayor, R. S. Montero, I. M. Llorente, and I. T. Foster, “Virtual infrastructure management in private and hybrid clouds,” *IEEE Internet Computing*, vol. 13, no. 5, pp. 14–22, 2009.
- [2]. G. Ateniese, R. C. Burns, R. Curtmola, J. Herring, L. Kissner, Z. N. J. Peterson, and D. X. Song, “Provable data possession at untrusted stores,” in *ACM Conference on Computer and Communications Security*, P. Ning, S. D. C. di Vimercati, and P. F. Syverson, Eds. ACM, 2007, pp. 598–609.
- [3]. A. Juels and B. S. K. Jr., “Pors: proofs of retrievability for large files,” in *ACM Conference on Computer and Communications Security*, P. Ning, S. D. C. di Vimercati, and P. F. Syverson, Eds. ACM, 2007, pp. 584–597.
- [4]. G. Ateniese, R. D. Pietro, L. V. Mancini, and G. Tsudik, “Scalable and efficient provable data possession,” in *Proceedings of the 4th international conference on Security and privacy in communication networks, SecureComm*, 2008, pp. 1–10.
- [5]. C. C. Erway, A. K. Upc, “u, C. Papamanthou, and R. Tamassia, “Dynamic provable data possession,” in *ACM Conference on Computer and Communications Security*, E. Al-Shaer, S. Jha, and A. D. Keromytis, Eds. ACM, 2009, pp. 213–222.
- [6]. H. Shacham and B. Waters, “Compact proofs of retrievability,” in *ASIACRYPT*, ser. Lecture Notes in Computer Science, J. Pieprzyk, Ed., vol. 5350. Springer, 2008, pp. 90–107.
- [7]. Q. Wang, C. Wang, J. Li, K. Ren, and W. Lou, “Enabling public verifiability and data dynamics for storage security in cloud computing,” in *ESORICS*, ser. Lecture Notes in Computer Science, M. Backes and P. Ning, Eds., vol. 5789. Springer, 2009, pp. 355–370.
- [8]. Y. Zhu, H. Wang, Z. Hu, G.-J. Ahn, H. Hu, and S. S. Yau, “Dynamic audit services for integrity verification of outsourced storages in clouds,” in *SAC*, W. C. Chu, W. E. Wong, M. J. Palakal, and C.-C. Hung, Eds. ACM, 2011, pp. 1550–1557.
- [9]. K. D. Bowers, A. Juels, and A. Oprea, “Hail: a high-availability and integrity layer for cloud storage,” in *ACM Conference on Computer and Communications Security*, E. Al-Shaer, S. Jha, and A. D. Keromytis, Eds. ACM, 2009, pp. 187–198.
- [10]. Y. Dodis, S. P. Vadhan, and D. Wichs, “Proofs of retrievability via hardness amplification,” in *TCC*, ser. Lecture Notes in Computer Science, O. Reingold, Ed., vol. 5444. Springer, 2009, pp. 109–127.
- [11]. L. Fortnow, J. Rompel, and M. Sipser, “On the power of multiprover interactive protocols,” in *Theoretical Computer Science*, 1988, pp. 156–161.
- [12]. Y. Zhu, H. Hu, G.-J. Ahn, Y. Han, and S. Chen, “Collaborative integrity verification in hybrid clouds,” in *IEEE Conference on the 7th International Conference on Collaborative Computing: Networking, Applications and Worksharing, CollaborateCom*, Orlando Florida, USA, October 15-18, 2011, pp. 197–206.
- [13]. M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. H. Katz, A. Konwinski, G. Lee, D. A. Patterson, A. Rabkin, I. Stoica, and M. Zaharia, “Above the clouds: A Berkeley view of cloud computing,” *EECS Department, University of California, Berkeley, Tech. Rep.*, Feb 2009.
- [14]. D. Boneh and M. Franklin, “Identity-based encryption from the weil pairing,” in *Advances in Cryptology (CRYPTO’2001)*, vol. 2139 of LNCS, 2001, pp. 213–229.

- [15]. O. Goldreich, Foundations of Cryptography: Basic Tools. Cambridge University Press, 2001.
- [16]. P. S. L. M. Barreto, S. D. Galbraith, C. O'Eigeartaigh, and M. Scott, "Efficient pairing computation on supersingular abelian varieties," Des. Codes Cryptography, vol. 42, no. 3, pp. 239–271, 2007.
- [17]. J.-L. Beuchat, N. Brisebarre, J. Detrey, and E. Okamoto, "Arithmetic operators for pairing-based cryptography," in CHES, ser. Lecture Notes in Computer Science, P. Paillier and I. Verbauwhede, Eds., vol. 4727. Springer, 2007, pp. 239–255.
- [18]. H. Hu, L. Hu, and D. Feng, "On a class of pseudorandom sequences from elliptic curves over finite fields," IEEE Transactions on Information Theory, vol. 53, no. 7, pp. 2598–2605, 2007.
- [19]. A. Bialecki, M. Cafarella, D. Cutting, and O. O'Malley, "Hadoop: A framework for running applications on large clusters built of commodity hardware," Tech. Rep., 2005. [Online]. Available: <http://lucene.apache.org/hadoop/>
- [20]. [20] E. Al-Shaer, S. Jha, and A. D. Keromytis, Eds., Proceedings of the 2009 ACM Conference on Computer and Communications Security, CCS 2009, Chicago, Illinois, USA, November 9-13, 2009. ACM, 2009.

About The Authors



B.SHANMUKHI, received his B.Tech degree in Information Technology from JNTU, Anantapur, India, in 2010. Currently pursuing M.Tech in computer science and engineering at Dr.KVSRCEW Institute of Technology, Kurnool, India.



D.SATYANARAYANA, received his M.Tech in Computer Science in MISTE from Jawaharlal Nehru Technological University, Anantapur, India, in 2011. He is an Asst.Professor at DR.K.V.S.R.C.E.W, Kurnool, India